

AND-Compression of NP-Complete Problems: Streamlined Proof and Minor Observations

Holger Dell^{1,2,3,4} 

Received: 23 September 2014 / Accepted: 23 December 2015
© Springer Science+Business Media New York 2016

Abstract Drucker (Proceedings of the 53rd annual symposium on foundations of computer science (FOCS), pp 609–618. doi:[10.1109/FOCS.2012.71](https://doi.org/10.1109/FOCS.2012.71), 2012) proved the following result: Unless the unlikely complexity-theoretic collapse $\text{coNP} \subseteq \text{NP/poly}$ occurs, there is no *AND-compression* for SAT. The result has implications for the compressibility and kernelizability of a whole range of NP-complete parameterized problems. We present a streamlined proof of Drucker’s theorem. An AND-compression is a deterministic polynomial-time algorithm that maps a set of SAT-instances x_1, \dots, x_t to a single SAT-instance y of size $\text{poly}(\max_i |x_i|)$ such that y is satisfiable if and only if all x_i are satisfiable. The “AND” in the name stems from the fact that the predicate “ y is satisfiable” can be written as the AND of all predicates “ x_i is satisfiable”. Drucker’s theorem complements the result by Bodlaender et al. (J Comput Syst Sci 75:423–434, 2009) and Fortnow and Santhanam (J Comput Syst Sci 77:91–106, 2011), who proved the analogous statement for OR-compressions, and Drucker’s proof not only subsumes their result but also extends it to *randomized* compression algorithms that are allowed to have a certain probability of failure. Drucker (Proceedings of the 53rd annual symposium on foundations of computer science (FOCS), pp 609–618. doi:[10.1109/FOCS.2012.71](https://doi.org/10.1109/FOCS.2012.71), 2012) presented two proofs: The first uses information theory and the minimax theorem from game theory, and the second is an elementary, iterative proof that is not as general. In our proof, we realize the iterative structure as a generalization of the arguments of Ko (J Comput Syst Sci

✉ Holger Dell
hdell@mmci.uni-saarland.de

¹ Cluster of Excellence (MMCI), Saarland University, Saarbrücken, Germany

² Simons Institute for the Theory of Computing, Berkeley, CA, USA

³ UC Berkeley, Berkeley, CA, USA

⁴ LIAFA, Université Paris Diderot, Paris, France

26:209–211, 1983) for P -selective sets, which use the fact that tournaments have dominating sets of logarithmic size. We generalize this fact to hypergraph tournaments. Our proof achieves the full generality of Drucker’s theorem, avoids the minimax theorem, and restricts the use of information theory to a single, intuitive lemma about the average noise sensitivity of compressive maps. To prove this lemma, we use the same information-theoretic inequalities as Drucker.

1 Introduction

Harnik and Naor [5] studied OR-compressions of SAT, and found cryptographic applications should they exist. The influential “OR-conjecture” by Bodlaender et al. [2] asserts they do not exist; more specifically, it asserts that t instances x_1, \dots, x_t of SAT cannot be mapped in polynomial time to a single instance y of size $\text{poly}(\max_i |x_i|)$ so that y is a yes-instance if and only if at least one x_i is a yes-instance. Assuming the truth of the OR-conjecture, the “composition framework” of Bodlaender et al. [2] has been used to show that many different problems in parameterized complexity do not have polynomial kernels. Fortnow and Santhanam [3] were able to prove that the OR-conjecture holds unless $\text{coNP} \subseteq \text{NP/poly}$, thereby connecting the OR-conjecture with a standard hypothesis in complexity theory.

The results of Bodlaender et al. [2] and Fortnow and Santhanam [3] can be used not only to rule out deterministic kernelization algorithms, but also to rule out randomized kernelization algorithms with one-sided error—this works as long as the success probability is bigger than zero even if it is exponentially small; thus, this is the same as allowing the kernelization algorithm to be a coNP -algorithm. It was left open whether the complexity-theoretic hypothesis $\text{coNP} \not\subseteq \text{NP/poly}$ (or some other hypothesis believed by complexity theorists) suffices to rule out kernelization algorithms that are randomized and have two-sided error. Drucker [1] resolves this question affirmatively; consequently, since almost all known polynomial kernel lower bounds can be based on Drucker’s theorem, we should now believe that none of these problems has a polynomial kernelization algorithm with a constant gap in its error probability.

With the same proof, Drucker [1] resolves a second important question: whether the “AND-conjecture”, which has also been formulated by Bodlaender et al. [2] analogous to the OR-conjecture, can be derived from existing complexity-theoretic assumptions. This is an intriguing question in itself, and it is also relevant for parameterized complexity as, for some parameterized problems, we can rule out polynomial kernels under the AND-conjecture, but we do not know how to do so under the OR-conjecture. Drucker [1] proves that the AND-conjecture is true if $\text{coNP} \not\subseteq \text{NP/poly}$ holds.

The purpose of this paper is to discuss Drucker’s theorem and its proof. To this end, we attempt to present a simpler proof of his theorem. Our proof in Sect. 3 gains in simplicity with a small loss in generality: the bound that we get is worse than Drucker’s bound by a factor of two. Using the slightly more complicated approach in Sect. 4, it is possible to get the same bounds as Drucker. These differences, however, do not matter for the basic version of the main theorem, which we state in Sect. 1.1

and further discuss in Sect. 1.2. For completeness, we briefly discuss a formulation of the composition framework in Sect. 1.3.

1.1 Main Theorem: Ruling Out OR- and AND-Compressions

An *OR-reduction* A for a language L is a polynomial-time reduction that maps a set $\{x_1, \dots, x_t\}$ to some instance $y \doteq A(\{x_1, \dots, x_t\})$ of a language L' such that y is a yes-instance of L' if and only if x_i is a yes-instance of L for some $i \in \{1, \dots, t\}$. An *AND-compression of size s* is an AND-reduction with bounded output size $|y| \leq s$. The terms *OR-reduction* and *OR-compression of size s* are defined analogously.

Every language has a trivial AND-compression of size $t \cdot n$, where n is defined as $\max_i |x_i|$. Drucker's theorem implies that we can only hope to improve this bound significantly for languages that are contained in $\text{NP/poly} \cap \text{coNP/poly}$. For example, any L with an AND-compression of size $t^{0.99} \cdot n^{0.99}$ is contained in $\text{NP/poly} \cap \text{coNP/poly}$. Previously, Fortnow and Santhanam [3] were only able to obtain the consequence that L is in NP/poly or in coNP/poly , neither of which is known to be closed under complement. In contrast, the complexity class $\text{NP/poly} \cap \text{coNP/poly}$ is closed under complement, which allows us to restrict our attention to OR-compressions for the remainder of this paper. To see that this is without loss of generality, note that every AND-compression for L is also an OR-compression for $\bar{L} \doteq \{0, 1\}^* \setminus L$ due to De Morgan's law: $y \in \bar{L}'$ holds if and only if $x_1 \in \bar{L}$ or $x_2 \in \bar{L}$ or \dots or $x_t \in \bar{L}$. Then proving that $\text{NP/poly} \cap \text{coNP/poly}$ contains \bar{L} is equivalent to proving that it contains L .

We now formally state Drucker's theorem.

Theorem 1 (Drucker's theorem) *Let $L, L' \subseteq \{0, 1\}^*$ be languages, let $\epsilon > 0$, and let $e_s, e_c \in [0, 1]$ with $e_s + e_c < 1$ be error probabilities. Assume that there exists a randomized polynomial-time algorithm A that maps any set $x = \{x_1, \dots, x_t\} \subseteq \{0, 1\}^n$ for some n and t to $y = A(x)$ such that:*

- (Soundness) *If all x_i 's are no-instances of L , then y is a no-instance of L' with probability $\geq 1 - e_s$.*
- (Completeness) *If exactly one x_i is a yes-instance of L , then y is a yes-instance of L' with probability $\geq 1 - e_c$.*
- (Size bound) *The size of y is bounded by $t^{1-\epsilon} \cdot \text{poly}(n)$.*

Then $L \in \text{NP/poly} \cap \text{coNP/poly}$.

The procedure A above does not need to be a “full” OR-compression, which makes the theorem more general. In particular, A is *relaxed* in two ways: it only needs to work, or be analyzed, in the case that all input instances have the same length; this is useful in hardness of kernelization proofs as it allows similar instances to be grouped together. Furthermore, A only needs to work, or be analyzed, in the case that at most one of the input instances is a yes-instance of L ; in an ongoing collaboration with Dániel Marx, we exploit this property to analyze the kernelization complexity of counting problems.

The fact that “relaxed” OR-compressions suffice in Theorem 1 is implicit in the proof of Drucker [1], but not stated explicitly. Before Drucker's work, Fortnow and

Santhanam [3] proved the special case of Theorem 1 in which $e_c = 0$, but they only obtain the weaker consequence $L \in \text{coNP/poly}$. Moreover, their proof uses the full completeness requirement and does not seem to work for relaxed OR-compressions.

1.2 Comparison and Overview of the Proof

The simplification of our proof stems from two main sources: (1) The “scaffolding” of our proof, its overall structure, is more modular and more similar to arguments used previously by Ko [4], Fortnow and Santhanam [3], and Dell and van Melkebeek [6] for compression-type procedures and Dell et al. [7] for isolation procedures. (2) While the information-theoretic part of our proof uses the same set of information-theoretic inequalities as Drucker’s, the simple version in Sect. 3 applies these inequalities to distributions that have a simpler structure. Moreover, our calculations have a somewhat more mechanical nature.

Both Drucker’s proof and ours use the relaxed OR-compression A to design a P/poly -reduction from L to the *statistical distance problem*, which is known to be in the intersection of NP/poly and coNP/poly by previous work (cf. Xiao [8]). Drucker [1] uses the minimax theorem and a game-theoretic sparsification argument to construct the polynomial advice of the reduction. He also presents an alternative proof [9, Section 3] in which the advice is constructed without these arguments and also without any explicit invocation of information theory; however, the alternative proof does not achieve the full generality of his theorem, and we feel that avoiding information theory entirely leads to a less intuitive proof structure. In contrast, our proof achieves full generality up to a factor of two in the simplest proof, it avoids game theoretic arguments, and it limits information theory to a single, intuitive lemma about the average noise sensitivity of compressive maps.

Using this information-theoretic lemma as a black box, we design the P/poly -reduction in a purely combinatorial way: We generalize the fact that tournaments have dominating sets of logarithmic size to *hypergraph tournaments*; these are complete t -uniform hypergraphs with the additional property that, for each hyperedge, one of its elements gets “selected”. In particular, for each set $e \subseteq \overline{L}$ of t no-instances, we select one element of e based on the fact that A ’s behavior on e somehow proves that the selected instance is a no-instance of L . The advice of the reduction is going to be a small dominating set of this hypergraph tournament on the set of no-instances of L . The crux is that we can efficiently test, with the help of the statistical distance problem oracle, whether an instance is dominated or not. Since any instance is dominated if and only if it is a no-instance of L , this suffices to solve L .

In the information-theoretic lemma, we generalize the notion of average noise sensitivity of Boolean functions (which can attain two values) to compressive maps (which can attain only relatively few values compared to the input length). We show that compressive maps have small average noise sensitivity. Drucker’s “distributional stability” is a closely related notion, which we make implicit use of in our proof. Using the latter notion as the anchor of the overall reduction, however, leads to some additional technicalities in Drucker’s proof, which we also run into in Sect. 4 where we obtain the same bounds as Drucker’s theorem. In Sect. 3 we instead use the average

noise sensitivity as the anchor of the reduction, which avoids these technicalities at the cost of losing a factor of two in the bounds.

1.3 Application: The Composition Framework for Ruling Out $O(k^{d-\epsilon})$ Kernels

We briefly describe a modern variant of the composition framework that is sufficient to rule out compressions and kernels of size $O(k^{d-\epsilon})$ using Theorem 1. Here, a problem is said to have *compressions of size s* if there exists a polynomial-time reduction to some problem such that the number of bits in the output of the reduction is at most s , and such a compression is called a *kernelization* if the reduction maps to instances of the same problem. Our framework for ruling out such compressions is almost identical to Lemma 1 of Dell and van Melkebeek [6] and Dell and Marx [10], and to Definition 2.2 of Hermelin and Wu [11]. By applying the framework for unbounded d , we can also use it to rule out polynomial kernels.

Definition 1 Let L be a language, let Π with parameter k be a parameterized problem, and let $d \geq 1$. A d -partite composition of L into Π is a polynomial-time algorithm A that maps any set $x = \{x_1, \dots, x_t\} \subseteq \{0, 1\}^n$ for some n and t to $y = A(x)$ such that:

- (1) If all x_i 's are no-instances of L , then y is a no-instance of Π .
- (2) If exactly one x_i is a yes-instance of L , then y is a yes-instance of Π .
- (3) The parameter k of y is bounded by $t^{1/d+o(1)} \cdot \text{poly}(n)$.

This notion of composition has an advantage over previous notions of OR-composition: The algorithm A does not need to work, or be analyzed, in the case that two or more of the x_i 's are yes-instances. The relaxation could make it easier in some cases to prove hardness of kernelization.

Definition 2 Let Π be a parameterized problem. We call Π d -compositional if there exists an NP-hard or coNP-hard problem L that has a d -partite composition algorithm into Π .

This definition encompasses both AND-compositions and OR-compositions because an AND-composition of L into Π is the same as an OR-composition of \bar{L} into $\bar{\Pi}$. Drucker [9, Theorems 7.7 and 7.11] applies his theorem to problems with polynomial kernels and to problems without them. Implicit in his arguments is the following corollary to Drucker's theorem.

Corollary 1 Let $d \geq 1$ and $\epsilon > 0$. If $\text{coNP} \not\subseteq \text{NP/poly}$, then no d -compositional problem has compressions or kernels of size $O(k^{d-\epsilon})$. Moreover, this even holds when the kernelization algorithm is allowed to be a randomized algorithm with two-sided, constant error probabilities e_s and e_c such that $e_s + e_c < 1$.

Proof We prove the contrapositive. Let L be an NP-hard or coNP-hard problem that has a d -partite composition A' into Π , and assume that Π has a kernelization algorithm with soundness error at most e_s and completeness error at most e_c so that $e_s + e_c < 1$ holds. We construct an OR-compression A that satisfies the conditions of Theorem 1. First, A runs the deterministic composition algorithm A' , which yields

an instance whose parameter k is bounded by $t^{1/d+o(1)} \cdot \text{poly}(n)$. On this instance, A runs the assumed randomized compression algorithm, which yields an instance of size $O(k^{d-\epsilon'})$ for some $\epsilon' > 0$. A simple calculation shows that this size is at most $t^{1-\epsilon} \cdot \text{poly}(n)$ for $\epsilon = \epsilon'/d$. By Theorem 1, we get $L \in (\text{coNP}/\text{poly} \cap \text{NP}/\text{poly})$ and thus $\text{coNP} \subseteq \text{NP}/\text{poly}$. \square

Several variants of the framework provided by this corollary are possible:

1. In order to rule out $\text{poly}(k)$ -kernels for a parameterized problem Π , we just need to prove that Π is d -compositional for all $d \in \mathbb{N}$; let's call Π *compositional* in this case. One way to show that Π is compositional is to construct a single *composition* from a hard problem L into Π ; this is an algorithm as in Definition 1, except that we replace (3) with the bound $k \leq t^{o(1)} \text{poly}(n)$.
2. Since all x_i 's in Definition 1 are promised to have the same length, we can consider a padded version \tilde{L} of the language L in order to filter the input instances of length n of the original L into a polynomial number of equivalence classes. Each input length of \tilde{L} in some interval $[p_1(n), p_2(n)]$ corresponds to one equivalence class of length- n instances of L . So long as \tilde{L} remains NP-hard or coNP-hard, it is sufficient to consider a composition from \tilde{L} into Π . Bodlaender et al. [12, Definition 4] formalize this approach.
3. The composition algorithm can also use randomness. In order to rule out deterministic kernelization algorithms under $\text{coNP} \not\subseteq \text{NP}/\text{poly}$, it is sufficient for the error probability $e_s + e_c$ of the composition algorithm to be bounded by a constant smaller than one.
4. In the case that L is NP-hard, Fortnow and Santhanam [3] and Dell and van Melkebeek [6] proved that the composition algorithm can also be a coNP-algorithm or even a coNP oracle communication game in order to get the collapse. Interestingly, this does not seem to follow from Drucker's proof nor from the proof presented here, and it seems to require the full completeness condition for the OR-composition. Kratsch [13] and Kratsch et al. [14] exploit these variants of the composition framework to prove kernel lower bounds.

2 Preliminaries

$A \dot{\cup} B$ refers to the union $A \cup B$ and stipulates that the sets A and B are disjoint. To define objects, we often write $\dot{=}$. For any set $R \subseteq \{0, 1\}^*$ and any $\ell \in \mathbb{N}$, we write $R_\ell \dot{=} R \cap \{0, 1\}^\ell$ for the set of all length- ℓ strings inside of R . For any $t \in \mathbb{N}$, we write $[t] \dot{=} \{1, \dots, t\}$. For a set V , we write $\binom{V}{\leq t}$ for the set of all subsets $x \subseteq V$ that have size at most t . We will work over a finite alphabet, usually $\Sigma = \{0, 1\}$. For a vector $a \in \Sigma^t$, a number $j \in [t]$, and a value $y \in \Sigma$, we write $a|_{j \leftarrow y}$ for the string that coincides with a except in position j , where it has value y . A *promise problem* is a tuple (Yes, No) with $\text{Yes}, \text{No} \subseteq \{0, 1\}^*$ and $\text{Yes} \cap \text{No} = \emptyset$; the elements of Yes are the *yes-instances* and the elements of No are the *no-instances* of the problem. A Boolean circuit $C : \{0, 1\}^n \rightarrow \{0, 1\}$ consists of n input gates, one output gate, AND-gates, and NOT-gates; the fan-in and depth of the circuits we consider in this paper are not restricted. For further background in complexity theory, we defer to the book

by Arora and Barak [15]. We assume some familiarity with the complexity classes **NP** and **coNP** as well as their non-uniform versions **NP/poly** and **coNP/poly**. In particular, we assume the following lemma, which is a popular homework assignment in complexity courses.

Lemma 1 *Let $L \subseteq \{0, 1\}^*$ be a language. Let*

$$\begin{aligned}\text{AND}(L) &= \{ (x_1, \dots, x_t) \mid \forall i \in [t]. x_i \in L \} , \text{ and} \\ \text{OR}(L) &= \{ (x_1, \dots, x_t) \mid \exists i \in [t]. x_i \in L \} .\end{aligned}$$

If $L \in \text{NP}$, then $\text{AND}(L) \in \text{NP}$ and $\text{OR}(L) \in \text{NP}$.

2.1 Distributions and Randomized Mappings

A *distribution* on a finite ground set Ω is a function $\mathcal{D}: \Omega \rightarrow [0, 1]$ with $\sum_{\omega \in \Omega} \mathcal{D}(\omega) = 1$. The *support* of \mathcal{D} is the set $\text{supp } \mathcal{D} = \{ \omega \in \Omega \mid \mathcal{D}(\omega) > 0 \}$. The *uniform distribution* \mathcal{U}_Ω on Ω is the distribution with $\mathcal{U}_\Omega(\omega) = \frac{1}{|\Omega|}$ for all $\omega \in \Omega$. We often view distributions as *random variables*, that is, we may write $f(\mathcal{D})$ to denote the distribution \mathcal{D}' that first produces a sample $\omega \sim \mathcal{D}$ and then outputs $f(\omega)$, where $f: \Omega \rightarrow \Omega'$. We use any of the following notations:

$$\mathcal{D}'(\omega') = \Pr(f(\mathcal{D}) = \omega') = \Pr_{\omega \sim \mathcal{D}}(f(\omega) = \omega') = \sum_{\omega \in \Omega} \mathcal{D}(\omega) \cdot \Pr(f(\omega) = \omega') .$$

The last term $\Pr(f(\omega) = \omega')$ in this equation is either 0 or 1 if f is a deterministic function, but we will also allow f to be a *randomized mapping*, that is, f has access to some “internal” randomness. This is modeled as a function $f: \Omega \times \{0, 1\}^r \rightarrow \Omega'$ for some $r \in \mathbb{N}$, and we write $f(\mathcal{D})$ as a short-hand for $f(\mathcal{D}, \mathcal{U}_{\{0, 1\}^r})$. That is, the internal randomness consists of a sequence of independent and fair coin flips.

2.2 Statistical Distance

The *statistical distance* $d(X, Y)$ between two distributions X and Y on Ω is defined as

$$d(X, Y) = \max_{T \subseteq \Omega} \left| \Pr(X \in T) - \Pr(Y \in T) \right| . \quad (1)$$

The statistical distance between X and Y is a number in $[0, 1]$, with $d(X, Y) = 0$ if and only if $X = Y$ and $d(X, Y) = 1$ if and only if the support of X is disjoint from the support of Y . It is an exercise to show the standard equivalence between the statistical distance and the 1-norm:

$$d(X, Y) = \frac{1}{2} \cdot \|X - Y\|_1 = \frac{1}{2} \sum_{\omega \in \Omega} \left| \Pr(X = \omega) - \Pr(Y = \omega) \right| .$$

2.3 The Statistical Distance Problem

The *statistical distance problem* $SD_{\leq \delta}^{\geq \Delta}$ is given two Boolean circuits $C, C' : \{0, 1\}^n \rightarrow \{0, 1\}$ to determine the distance $d(C(\mathcal{U}_{\{0,1\}^n}), C'(\mathcal{U}_{\{0,1\}^n}))$ between the output distributions of the circuits. More formally, let $\delta, \Delta : \mathbb{N} \rightarrow [0, 1]$ be functions so that $\delta(n) < \Delta(n)$ holds for all $n \in \mathbb{N}$. Then $SD_{\leq \delta}^{\geq \Delta}$ is defined as the promise problem whose instances are pairs (C, C') of Boolean circuits $C, C' : \{0, 1\}^n \rightarrow \{0, 1\}^*$; the yes-instances are pairs (C, C') with a relatively different output distribution, that is, pairs that satisfy $d(C(\mathcal{U}_{\{0,1\}^n}), C'(\mathcal{U}_{\{0,1\}^n})) \geq \Delta(n)$; the no-instances are pairs with a relatively similar output distribution, that is, pairs that satisfy $d(C(\mathcal{U}_{\{0,1\}^n}), C'(\mathcal{U}_{\{0,1\}^n})) \leq \delta(n)$.

The statistical distance problem is not known to be polynomial-time computable, and in fact it is not believed to be. On the other hand, the problem is also not believed to be NP-hard because the problem is computationally easy in the following sense.

Theorem 2 (Xiao [8] + Adleman [16])

If δ and Δ are constants with $0 \leq \delta < \Delta \leq 1$, then $SD_{\leq \delta}^{\geq \Delta} \in (\text{NP/poly} \cap \text{coNP/poly})$.

Moreover, the same holds when $\delta(n)$ and $\Delta(n)$ are functions with $\Delta - \delta \geq \frac{1}{\text{poly}(n)}$.

Note that $SD_{\leq \delta}^{\geq \Delta}$ is a promise problem and NP is usually defined as a class of languages. To ease the notation in Theorem 2, we write NP for the class of all promise problems $P = (\text{Yes}, \text{No})$ such that there is a polynomial-time many-one reduction from P to the satisfiability problem; similarly, we extend coNP to promise problems.

Slightly stronger versions of Theorem 2 are known: For example, Xiao [8, p. 144ff] proves that $SD_{\leq \delta}^{\geq \Delta} \in \text{AM} \cap \text{coAM}$ holds. In fact, Theorem 2 is established by combining his theorem with the standard fact that $\text{AM} \subseteq \text{NP/poly}$, i.e., that Arthur–Merlin games can be derandomized with polynomial advice [16]. Moreover, when we have the stronger guarantee $\Delta^2 - \delta > \frac{1}{\text{poly}(n)}$, then $SD_{\leq \delta}^{\geq \Delta}$ can be solved using *statistical zero-knowledge proof systems* [17, 18]. Finally, if $\Delta = 1 > \delta + \frac{1}{\text{poly}(n)}$, the problem can be solved with *perfect zero-knowledge proof systems* [17, Proposition 5.7]. Using these stronger results whenever possible gives slightly stronger complexity collapses in the main theorem.

3 Ruling Out OR-Compressions

In this section we prove Theorem 1: Any language L that has a relaxed OR-compression is in $\text{coNP/poly} \cap \text{NP/poly}$. We rephrase the theorem in a form that reveals the precise inequality between the error probabilities and the compression ratio needed to get the complexity consequence.

Theorem 3 (ϵt -compressive version of Drucker’s theorem) *Let $L, L' \subseteq \{0, 1\}^*$ be languages and $e_s, e_c \in [0, 1]$ be some constants denoting the error probabilities. Let $t = t(n) > 0$ be a polynomial and $\epsilon > 0$. Let*

$$A: \left(\begin{smallmatrix} \{0, 1\}^n \\ \leq t \end{smallmatrix} \right) \rightarrow \{0, 1\}^{\epsilon t} \quad (2)$$

be a randomized $\mathbf{P/poly}$ -algorithm such that, for all $x \in \left(\begin{smallmatrix} \{0, 1\}^n \\ \leq t \end{smallmatrix} \right)$,

- if $|x \cap L| = 0$, then $A(x) \in \overline{L'}$ holds with probability $\geq 1 - e_s$, and
- if $|x \cap L| = 1$, then $A(x) \in L'$ holds with probability $\geq 1 - e_c$.

If $e_s + e_c < 1 - \sqrt{(2 \ln 2)\epsilon}$, then $L \in \mathbf{NP/poly} \cap \mathbf{coNP/poly}$.

This is Theorem 7.1 in Drucker [9]. However, there are two noteworthy differences:

1. Drucker obtains complexity consequences even when $e_s + e_c < 1 - \sqrt{(\ln 2/2)\epsilon}$ holds, which makes his theorem more general. The difference stems from the fact that we optimized the proof in this section for simplicity and not for the optimality of the bound. He also obtains complexity consequences under the (incomparable) bound $e_s + e_c < 2^{-\epsilon-3}$. Using the slightly more complicated setup of Sect. 4, we would be able to achieve both of these bounds.
2. To get a meaningful result for OR-compression of \mathbf{NP} -complete problems, we need the complexity consequence $L \in \mathbf{NP/poly} \cap \mathbf{coNP/poly}$ rather than just $L \in \mathbf{NP/poly}$. To get the stronger consequence, Drucker relies on the fact that the statistical distance problem $\text{SD}_{\leq \delta}^{\geq \Delta}$ has statistical zero knowledge proofs. If Δ and δ are constants, this is only known to be true when $\Delta^2 > \delta$ holds, which translates to the more restrictive assumption $(e_s + e_c)^2 < 1 - \sqrt{(\ln 2/2)\epsilon}$ in his theorem. We instead use Theorem 2, which does not go through statistical zero knowledge and proves more directly that $\text{SD}_{\leq \delta}^{\geq \Delta}$ is in $\mathbf{NP/poly} \cap \mathbf{coNP/poly}$ whenever $\Delta > \delta$ holds. Doing so in Drucker's paper immediately improves all of his $L \in \mathbf{NP/poly}$ consequences to $L \in \mathbf{NP/poly} \cap \mathbf{coNP/poly}$.

To obtain Theorem 1, the basic version of Drucker's theorem, as a corollary of Theorem 3, none of these differences matter. This is because we could choose $\epsilon > 0$ to be sufficiently smaller in the proof of Theorem 1, which we provide now before we turn to the proof of Theorem 3.

Proof (of Theorem 1) Let A be the algorithm assumed in Theorem 1, and let $C \geq 2$ be large enough so that the output size of A is bounded by $t^{1-1/C} \cdot C \cdot n^C$. We transform A into an algorithm as required for Theorem 3. Let $\epsilon > 0$ be a small enough constant so that $e_s + e_c < 1 - \sqrt{(2 \ln 2)\epsilon}$. Moreover, let $t(n)$ be a large enough polynomial so that $(t(n))^{1-1/C} \cdot C \cdot n^C < \epsilon t(n)$ holds. Then we restrict A to a family of functions $A_n: \left(\begin{smallmatrix} \{0, 1\}^n \\ \leq t(n) \end{smallmatrix} \right) \rightarrow \{0, 1\}^{< \epsilon t(n)}$. Each of these functions can be computed by polynomial-size circuits with additional random bits, and thus the family $(A_n)_{n \in \mathbb{N}}$ is a randomized $\mathbf{P/poly}$ -algorithm. Now a minor observation is needed to get an algorithm of the form (2): The set $\{0, 1\}^{< \epsilon t}$ can be efficiently encoded in $\{0, 1\}^{\epsilon t}$, which changes the output language from L' to some L'' . Overall, we constructed a family A_n that satisfies the assumptions of Theorem 3, and we obtain $L \in \mathbf{NP/poly} \cap \mathbf{coNP/poly}$, which proves the claim of Theorem 1. \square

3.1 ORs are Sensitive to Yes-Instances

The semantic property of relaxed OR-compressions is that they are “ L -sensitive”: They show a dramatically different behavior for all-no input sets versus input sets that contain a single yes-instance of L . The following simple fact is the only place in the overall proof where we use the soundness and completeness properties of A .

Lemma 2 *For all distributions X on $\binom{\bar{L}}{<t}$ and all $v \in L$, we have*

$$d\left(A(X), A(X \cup \{v\})\right) \geq \Delta \doteq 1 - (e_s + e_c). \quad (3)$$

Proof The probability that $A(X)$ outputs an element of L' is at most e_s , and similarly, the probability that $A(X \cup \{v\})$ outputs an element of L' is at least $1 - e_c$. By (1) with $T = L'$, the statistical distance between the two distributions is at least Δ . \square

Despite the fact that relaxed OR-compressions are sensitive to the presence or absence of a yes-instance, we argue next that their behavior *within* the set of no-instances is actually quite predictable.

3.2 The Average Noise Sensitivity of Compressive Maps is Small

Relaxed OR-compressions are in particular compressive maps. The following lemma says that the average noise sensitivity of any compressive map is low. Here, “average noise sensitivity” refers to the difference in the behavior of a function when the input is subject to random noise; in our case, we change the input in a single random location and notice that the behavior of a compressive map does not change much.

Lemma 3 *Let t be a positive integer, let X be the uniform distribution on $\{0, 1\}^t$, and let $\epsilon > 0$. Then, for all randomized mappings $f: \{0, 1\}^t \rightarrow \{0, 1\}^{\epsilon t}$, we have*

$$\mathbf{E}_{j \sim \mathcal{U}_{[t]}} d\left(f(X|_{j \leftarrow 0}), f(X|_{j \leftarrow 1})\right) \leq \delta \doteq \sqrt{2 \ln 2 \cdot \epsilon}. \quad (4)$$

We defer the purely information-theoretic proof of this lemma to Sect. 3.5. In the special case where $f: \{0, 1\}^t \rightarrow \{0, 1\}$ is a Boolean function, the left-hand side of (4) coincides with the usual definition of the average noise sensitivity.

We translate Lemma 3 to our relaxed OR-compression A as follows.

Lemma 4 *Let $t, n \in \mathbb{N}$ with $0 < t \leq 2^n$, let $\epsilon > 0$, and let $A: \binom{\{0,1\}^n}{\leq t} \rightarrow \{0, 1\}^{\epsilon t}$ be any randomized mapping. For all $e \in \binom{\{0,1\}^n}{t}$, there exists $v \in e$ so that*

$$d\left(A(\mathcal{U}_{2^e} \setminus \{v\}), A(\mathcal{U}_{2^e} \cup \{v\})\right) \leq \delta. \quad (5)$$

Here \mathcal{U}_{2^e} samples a subset of e uniformly at random. Note that we replaced the expectation over j from (4) with the mere existence of an element v in (5) since this is all we need; the stronger property also holds.

Proof To prove the claim, let v_1, \dots, v_t be the elements of e in lexicographic order. For $b \in \{0, 1\}^t$, let $g(b) \subseteq e$ be such that $v_i \in g(b)$ holds if and only if $b_i = 1$. We define the randomized mapping $f: \{0, 1\}^t \rightarrow \{0, 1\}^{\epsilon t}$ as follows:

$$f(b_1, \dots, b_t) \doteq A(g(b)).$$

Then $f(X|_{j \leftarrow 0}) = A(\mathcal{U}_{2^e} \setminus \{v_j\})$ and $f(X|_{j \leftarrow 1}) = A(\mathcal{U}_{2^e} \cup \{v_j\})$. The claim follows from Lemma 3 with $v \doteq v_j$ for some j that minimizes the statistical distance in (4). \square

This lemma suggests the following tournament idea. We let $V = \overline{L}_n$ be the set of no-instances, and we let them compete in matches consisting of t players each. That is, a match corresponds to a hyperedge $e \in \binom{V}{t}$ of size t and every such hyperedge is present, so we are looking at a complete t -uniform hypergraph. We say that a player $v \in e$ is “selected” in the hyperedge e if the behavior of A on $\mathcal{U}_{2^e} \setminus \{v\}$ is not very different from the behavior of A on $\mathcal{U}_{2^e} \cup \{v\}$, that is, if (5) holds. The point of this construction is that v being selected proves that v must be a no-instance because (3) does not hold. We obtain a “selector” function $S: \binom{V}{t} \rightarrow V$ that, given e , selects an element $v = S(e) \in e$. We call S a *hypergraph tournament* on V .

3.3 Hypergraph Tournaments Have Small Dominating Sets

Tournaments are complete directed graphs, and it is well-known that they have dominating sets of logarithmic size (e.g., see [4]). A straightforward generalization applies to hypergraph tournaments $S: \binom{V}{t} \rightarrow V$. We say that a set $g \in \binom{V}{t-1}$ *dominates* a vertex v if $v \in g$ or $S(g \cup \{v\}) = v$ holds. A set $\mathcal{D} \subseteq \binom{V}{t-1}$ is a *dominating set* of S if all vertices $v \in V$ are dominated by at least one element in \mathcal{D} .

Lemma 5 *Let V be a finite set, and let $S: \binom{V}{t} \rightarrow V$ for some integer $t \geq 2$ be a hypergraph tournament. Then S has a dominating set $\mathcal{D} \subseteq \binom{V}{t-1}$ of size at most $t \log |V|$.*

Proof We construct the set \mathcal{D} inductively. Initially, it has $k = 0$ elements. After the k -th step of the construction, we will preserve the invariant that \mathcal{D} is of size exactly k and that $|R| \leq (1 - 1/t)^k \cdot |V|$ holds, where R is the set of vertices that are not yet dominated, that is,

$$R = \{v \in V \mid v \notin g \text{ and } S(g \cup \{v\}) \neq v \text{ holds for all } g \in \mathcal{D}\}.$$

If $0 < |R| < t$, we can add an arbitrary edge $g^* \in \binom{V}{t-1}$ with $R \subseteq g^*$ to \mathcal{D} to finish the construction. Otherwise, the following averaging argument, shows that there is an element $g^* \in \binom{R}{t-1}$ that dominates at least a $1/t$ -fraction of elements $v \in R$:

$$\frac{1}{t} = \mathbf{E}_{e \in \binom{R}{t}} \Pr_{v \in e} (S(e) = v) = \mathbf{E}_{g \in \binom{R}{t-1}} \Pr_{v \in R-g} (S(g \cup \{v\}) = v).$$

Thus, the number of elements of R left undominated by g^* is at most $(1 - 1/t) \cdot |R|$, so the inductive invariant holds after including g^* into \mathcal{D} . Since $(1 - 1/t)^k \cdot |V| \leq \exp(-k/t) \cdot |V| < 1$ for $k = t \log |V|$, we have $R = \emptyset$ after $k \leq t \log |V|$ steps of the construction, and, in particular, \mathcal{D} has at most $t \log |V|$ elements. \square

3.4 Proof of the Main Theorem: Reduction to Statistical Distance

Proof (of Theorem 3) We describe a deterministic P/poly reduction from L to the statistical distance problem $\text{SD}_{\leq \delta}^{\geq \Delta}$ with $\Delta = 1 - (e_s + e_c)$ and $\delta = \sqrt{(2 \ln 2)\epsilon}$. The reduction outputs the conjunction of polynomially many instances of $\text{SD}_{\leq \delta}^{\geq \Delta}$. By Theorem 2, $\text{SD}_{\leq \delta}^{\geq \Delta}$ polynomial-time many-one reduces to a problem in the intersection of NP/poly and coNP/poly. By Lemma 1, this intersection is closed under taking conjunctions, and we obtain $L \in \text{NP/poly} \cap \text{coNP/poly}$. Thus it remains to find such a reduction to the statistical distance problem. To simplify the discussion, we describe the reduction in terms of an algorithm that solves L and uses $\text{SD}_{\leq \delta}^{\geq \Delta}$ as an oracle. However, the algorithm only makes non-adaptive queries at the end of the computation and accepts if and only if all oracle queries accept; this corresponds to a reduction that maps an instance of L to a conjunction of instances of $\text{SD}_{\leq \delta}^{\geq \Delta}$ as required.

To construct the advice at input length n , we use Lemma 4 with $t = t(n)$ to obtain a hypergraph tournament S on $V = \overline{L}_n$, which in turn gives rise to a small dominating set $\mathcal{D} \subseteq \binom{V}{t-1}$ by Lemma 5. We remark that if $|V| \leq t = \text{poly}(n)$, then we can use V , the set of all no-instances of L at this input length, as the advice. Otherwise, we define the hypergraph tournament S for all $e \in \binom{V}{t}$ as follows:

$$S(e) \doteq \min \{ v \in e \mid d(A(\mathcal{U}_{2^e} \setminus \{v\}), A(\mathcal{U}_{2^e} \cup \{v\})) \leq \delta \}.$$

Here, the minimum returns the lexicographically first v in the set. By Lemma 4, the set over which the minimum is taken is non-empty, and thus S is well-defined. Furthermore, the hypergraph tournament has a dominating set \mathcal{D} of size at most tn by Lemma 5. As advice for input length n , we choose this set \mathcal{D} . Now we have $v \in \overline{L}$ if and only if v is dominated by \mathcal{D} . The idea of the reduction is to efficiently check the latter property.

The algorithm works as follows: Let $v \in \{0, 1\}^n$ be an instance of L given as input. If $v \in g$ holds for some $g \in \mathcal{D}$, the algorithm rejects v and halts. Otherwise, it queries the SD-oracle on the instance $(A(\mathcal{U}_{2^g}), A(\mathcal{U}_{2^g} \cup \{v\}))$ for each $g \in \mathcal{D}$. If the oracle claims that all queries are yes-instances, our algorithm accepts, and otherwise, it rejects.

First note that distributions of the form $A(\mathcal{U}_{2^g})$ and $A(\mathcal{U}_{2^g} \cup \{v\})$ can be sampled by using polynomial-size circuits, and so they form syntactically correct instances of the SD-problem. The information about A , g , and v is hard-wired into these circuits. More precisely, we construct two circuit families, one for each distribution. Each circuit family is based on the circuit representation of the P/poly-algorithm A ; this yields a circuit family whose size grows at most polynomially. The input bits of the circuits are used to produce a sample from \mathcal{U}_{2^g} for the first distribution, and a sample

from $\mathcal{U}_{2g} \cup \{v\}$ for the second distribution. Further input bits play the role of the internal randomness of A in case A is a randomized algorithm.

It remains to prove the correctness of the reduction. If $v \in L$, we have for all $g \in \mathcal{D} \subseteq \bar{L}$ that $v \notin g$ and that the statistical distance of the query corresponding to g is at least $\Delta = 1 - (e_s + e_c)$ by Lemma 2. Thus all queries that the reduction makes satisfy the promise of the SD-problem and the oracle answers the queries correctly, leading our reduction to accept. On the other hand, if $v \notin L$, then, since \mathcal{D} is a dominating set of \bar{L} with respect to the hypergraph tournament S , there is at least one $g \in \mathcal{D}$ so that $v \in g$ or $S(g \cup \{v\}) = v$ holds. If $v \in g$, the reduction rejects. The other case implies that the statistical distance between $A(\mathcal{U}_{2g})$ and $A(\mathcal{U}_{2g} \cup \{v\})$ is at most δ . The query corresponding to this particular g therefore satisfies the promise of the SD-problem, which means that the oracle answers correctly on this query and our reduction rejects. \square

3.5 Information-Theoretic Arguments

We now prove Lemma 3. The proof uses the *Kullback–Leibler divergence* as an intermediate step. Just like the statistical distance, this notion measures how similar two distributions are, but it does so in an information-theoretic way rather than in a purely statistical way. In fact, it is well-known in the area that the Kullback–Leibler divergence and the mutual information are almost interchangeable in a certain sense. Drucker’s techniques yield a version of this paradigm, which we formally state in Lemma 6 below; then we prove Lemma 3 by bounding the statistical distance in terms of the Kullback–Leibler divergence using standard inequalities.

We introduce some basic information-theoretic notions. The *Shannon entropy* $H(X)$ of a random variable X is

$$H(X) = \mathbf{E}_{x \sim X} \log \left(\frac{1}{\Pr(X = x)} \right).$$

The conditional Shannon entropy $H(X|Y)$ is

$$\begin{aligned} H(X|Y) &= \mathbf{E}_{y \sim Y} H(X|Y = y) \\ &= \mathbf{E}_{y \sim Y} \sum_x \Pr(X = x \mid Y = y) \cdot \log \left(\frac{1}{\Pr(X = x \mid Y = y)} \right). \end{aligned}$$

The *mutual information* between X and Y is $I(X : Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$. Note that $I(X : Y) \leq \log |\text{supp } X|$, where $|\text{supp } X|$ is the size of the support of X . The *conditional mutual information* can be defined by the *chain rule of mutual information* $I(X : Y \mid Z) = I(X : YZ) - I(X : Z)$. If Y and Z are independent, then a simple calculation reveals that $I(X : Y) \leq I(X : Y \mid Z)$ holds.

Before we proceed, we further define the binary *Kullback–Leibler divergence*. Let X and Y be two distributions so that $\text{supp}(X) \subseteq \text{supp}(Y)$. Then

$$D_{\text{KL}}(X \parallel Y) \doteq \sum_{z \in \text{supp}(X)} \log \left(\frac{\Pr(X = z)}{\Pr(Y = z)} \right) \cdot \Pr(X = z) . \quad (6)$$

The KL-divergence can be interpreted as a kind of distance between two distributions, and in fact, Pinsker's inequality [19, Lemma 11.6.1] gives the following connection to the statistical distance between X and Y :

$$d(X, Y) \leq \sqrt{\frac{\ln 2}{2} \cdot D_{\text{KL}}(X \parallel Y)} . \quad (7)$$

Thus when the KL-divergence is smaller than $2/\ln 2$, Pinsker's inequality yields a non-trivial bound on the statistical distance.

We now establish a bound on the Kullback–Leibler divergence. The application of Lemma 3 only uses $\Sigma = \{0, 1\}$. The proof does not become more complicated for general Σ , and we will need the more general version later in this paper.

Lemma 6 *Let t be a positive integer and let X_1, \dots, X_t be independent distributions on some finite set Σ , and let $X = X_1, \dots, X_t$. Then, for all randomized mappings $f: \Sigma^t \rightarrow \{0, 1\}^*$, we have the following upper bound on the expected value of the Kullback–Leibler divergence:*

$$\mathbf{E}_{j \sim \mathcal{U}_{[t]}} \mathbf{E}_{x \sim X_j} D_{\text{KL}}(f(X|_{j \leftarrow x}) \parallel f(X)) \leq \frac{1}{t} \cdot I(f(X) : X) .$$

Proof The result follows by a basic calculation with entropy notions. The first equality is the definition of the Kullback–Leibler divergence, which we rewrite using the logarithm rule $\log(a/b) = \log(1/b) - \log(1/a)$ and the linearity of expectation:

$$\begin{aligned} & \mathbf{E}_{j \sim \mathcal{U}_{[t]}} \mathbf{E}_{x \sim X_j} D_{\text{KL}}(f(X|_{j \leftarrow x}) \parallel f(X)) \\ &= \mathbf{E}_j \mathbf{E}_x \sum_z \log \left(\frac{\Pr(f(X|_{j \leftarrow x}) = z)}{\Pr(f(X) = z)} \right) \cdot \Pr(f(X|_{j \leftarrow x}) = z) \\ &= \mathbf{E}_j \sum_z \log \left(\frac{1}{\Pr(f(X) = z)} \right) \cdot \mathbf{E}_x \Pr(f(X|_{j \leftarrow x}) = z) \\ &\quad - \mathbf{E}_j \mathbf{E}_x \sum_z \log \left(\frac{1}{\Pr(f(X|_{j \leftarrow x}) = z)} \right) \cdot \Pr(f(X|_{j \leftarrow x}) = z) . \end{aligned}$$

Now $\mathbf{E}_x \Pr(f(X|_{j \leftarrow x}) = z) = \Pr(f(X) = z)$ holds and, by independence of the X_j 's, we have $X|_{j \leftarrow x} = (X|X_j = x)$. Hence both terms of the sum above can be expressed as entropies:

$$\begin{aligned}
 \dots &= H(f(X)) - \mathbf{E}_{j \sim \mathcal{U}_{[t]}} \mathbf{E}_{x \sim X_j} H(f(X) | X_j = x) && \text{(definition of entropy)} \\
 &= H(f(X)) - \mathbf{E}_j H(f(X) | X_j) && \text{(definition of conditional entropy)} \\
 &= \mathbf{E}_j I(f(X) : X_j) && \text{(definition of mutual information)} \\
 &\leq \frac{1}{t} \cdot \sum_{j \in [t]} I(f(X) : X_j \mid X_1 \dots X_{j-1}) && \text{(by independence of } X'_j\text{'s)} \\
 &= \frac{1}{t} \cdot I(f(X) : X). && \text{(chain rule of mutual information)}
 \end{aligned}$$

□

We now turn to the proof of Lemma 3, where we bound the statistical distance in terms of the Kullback–Leibler divergence.

Proof (of Lemma 3) We observe that $I(f(X) : X) \leq \log |\text{supp } f(X)| \leq \epsilon t$, and so we are in the situation of Lemma 6 with $\Sigma = \{0, 1\}$. We first apply the triangle inequality to the left-hand side of (4). Then we use Pinsker’s inequality to bound the statistical distance in terms of the Kullback–Leibler divergence, which we can in turn bound by ϵ using Lemma 6.

$$\begin{aligned}
 &\mathbf{E}_{j \sim \mathcal{U}_{[t]}} d\left(f(X|_{j \leftarrow 0}), f(X|_{j \leftarrow 1})\right) \\
 &\leq \mathbf{E}_{j \sim \mathcal{U}_{[t]}} d\left(f(X), f(X|_{j \leftarrow 0})\right) && \text{(triangle inequality)} \\
 &\quad + \mathbf{E}_{j \sim \mathcal{U}_{[t]}} d\left(f(X), f(X|_{j \leftarrow 1})\right) \\
 &= 2 \cdot \mathbf{E}_{j \sim \mathcal{U}_{[t]}} \mathbf{E}_{x \sim X_j} d\left(f(X), f(X|_{j \leftarrow x})\right) \\
 &\leq 2 \cdot \mathbf{E}_j \mathbf{E}_x \sqrt{\frac{\ln 2}{2} \cdot D_{\text{KL}}\left(f(X|_{j \leftarrow x}) \parallel f(X)\right)} && \text{(Pinsker’s inequality)} \\
 &\leq 2 \cdot \sqrt{\frac{\ln 2}{2}} \cdot \mathbf{E}_j \mathbf{E}_x D_{\text{KL}}\left(f(X|_{j \leftarrow x}) \parallel f(X)\right) && \text{(Jensen’s inequality)} \\
 &\leq 2 \cdot \sqrt{\frac{\ln 2}{2}} \cdot \epsilon = \delta. && \text{(Lemma 6)}
 \end{aligned}$$

The equality above uses the fact that X_j is the uniform distribution on $\{0, 1\}$. □

4 Extension: Ruling Out OR-Compressions of Size $O(t \log t)$

In this section we tweak the proof of Theorem 3 so that it works even when the t instances of L are mapped to an instance of L' of size at most $O(t \log t)$. The drawback is that we cannot handle positive constant error probabilities for randomized

relaxed OR-compression anymore. For simplicity, we restrict ourselves to *deterministic* relaxed OR-compressions of size $O(t \log t)$ throughout this section.

Theorem 4 ($O(t \log t)$ -compressive version of Drucker's theorem)

Let $L, L' \subseteq \{0, 1\}^*$ be languages. Let $t = t(n) > 0$ be a polynomial. Assume there exists a P/poly-algorithm

$$A: \left(\begin{array}{c} \{0, 1\}^n \\ \leq t \end{array} \right) \rightarrow \{0, 1\}^{O(t \log t)}$$

such that, for all $x \in \left(\begin{array}{c} \{0, 1\}^n \\ \leq t \end{array} \right)$,

- if $|x \cap L| = 0$, then $A(x) \in \overline{L'}$, and
- if $|x \cap L| = 1$, then $A(x) \in L'$.

Then $L \in \text{NP/poly} \cap \text{coNP/poly}$.

This is Theorem 7.1 in Drucker [9]. The main reason why the proof in Sect. 3 breaks down for compressions to size ϵt with $\epsilon = O(\log t)$ is that the bound on the statistical distance in Lemma 3 becomes trivial. This happens already when $\epsilon \geq \frac{1}{2 \ln 2} \approx 0.72$. Nevertheless, the bound that Lemma 6 gives for the Kullback–Leibler divergence remains non-trivial even for relatively large $\epsilon = O(\log t)$.

The proof of Lemma 3 relies on Pinsker's inequality, which becomes trivial in the parameter range under consideration. For this reason, Drucker [9] uses a different inequality between the statistical distance and the Kullback–Leibler divergence, Vajda's inequality, that still gives a non-trivial bound on the statistical distance when the divergence is $\geq \frac{1}{2 \ln 2}$. The inequality works out such that if the divergence is logarithmic, then the statistical distance is an inverse polynomial away from 1. We obtain the following analogue to Lemma 3.

Lemma 7 Let $\Gamma > 0$. There exists a large enough positive integer c such that, for all positive integers t , all independent uniform distributions X_1, \dots, X_t on some finite set Σ and their joint distribution $X = X_1, \dots, X_t$, and, for all randomized mappings $f: \Sigma^t \rightarrow \{0, 1\}^*$ with $I(f(X) : X) \leq \Gamma \cdot t \log t$, we have

$$\mathbf{E}_{j \sim \mathcal{U}_{[t]}} \mathbf{E}_{a \sim X_j} d\left(f(X|_{X_j \neq a}), f(X|_{X_j = a})\right) \leq 1 - \frac{1}{ct^c} + \frac{1}{|\Sigma|}. \quad (8)$$

The notation $X|_{X_j \neq x}$ refers to the random variable that samples $x_i \sim X_i = \mathcal{U}_\Sigma$ independently for each $i \neq j$ as usual, and that samples x_j from the distribution X_j conditioned on the event that $X_j \neq a$, that is, the distribution $\mathcal{U}_{\Sigma \setminus \{a\}}$. The notation $X|_{X_j = a} = X|_{j \leftarrow a}$ is as before, that is, $x_j = a$ is fixed.

We defer the proof of the lemma to the end of this section and discuss now how to use it to obtain the stronger result for $O(t \log t)$ compressions. First note that we could not have directly used Lemma 7 in place of Lemma 3 in the proof of the main result, Theorem 3. This is because for $\Sigma = \{0, 1\}$, the right-hand side of (8) becomes bigger than 1 and thus trivial. In fact, this is the reason why we formulated Lemma 6 for general Σ . We need to choose Σ with $|\Sigma| = Ct^C$ for some large enough constant C to get anything meaningful out of (8).

4.1 A Different Hypergraph Tournament

To be able to work with larger Σ , we need to define the hypergraph tournament in a different way; not much is changing on a conceptual level, but the required notation becomes a bit less natural. We do this as follows.

Lemma 8 *Let $\Gamma > 0$. There exist a large enough positive integer C such that the following holds for all $t, n \in \mathbb{N}$ with $0 < t \leq 2^n$ and all randomized mappings $A: \binom{\{0,1\}^n}{\leq t} \rightarrow \{0,1\}^{\Gamma \cdot t \log t}$. For all $e = e_1 \dot{\cup} e_2 \dot{\cup} \dots \dot{\cup} e_t \subseteq \{0,1\}^n$ with $|e_1| = \dots = |e_t| = Ct^C$, there is an element $v \in e$ so that*

$$d\left(A(X_e|_{v \notin X_e}), A(X_e|_{v \in X_e})\right) \leq 1 - \frac{1}{Ct^C}, \quad (9)$$

where X_e is the distribution that samples the t -element set $\{\mathcal{U}_{e_1}, \dots, \mathcal{U}_{e_t}\}$, and $X_e|_E$ is the distribution X_e conditioned on the event E .

For instance if $v \in e_1$, then $X_e|_{v \notin X_e}$ samples the t -element set $\{\mathcal{U}_{e_1 \setminus \{v\}}, \mathcal{U}_{e_2}, \dots, \mathcal{U}_{e_t}\}$ and $X_e|_{v \in X_e}$ samples the t -element set $\{v, \mathcal{U}_{e_2}, \dots, \mathcal{U}_{e_t}\}$. The proof of this lemma is analogous to the proof of Lemma 4.

Proof We choose C large enough later. Let $\Sigma = [Ct^C]$. From A , we define a function f by restricting A to the support of X_e . More precisely, let $e_{ia} \in \{0,1\}^n$ for $i \in [t]$ and $a \in \Sigma$ be the lexicographically a -th element of e_i . We define the function $f: \Sigma^t \rightarrow \{0,1\}^{\Gamma t \log t}$ as follows: $f(a_1, \dots, a_t) \doteq A(e_{1a_1}, \dots, e_{ta_t})$. Then $I(f(X) : X) \leq \Gamma t \log t$ holds. Finally, for all $i \in [t]$, let the X_i be independent uniform distributions on Σ . We apply Lemma 7 to f and obtain (8). Let $1 - 1/(ct^c) + 1/|\Sigma|$ be the right-hand side of (8) applied to f . Here, c is just a function of Γ , and we set C such that $1 - 1/(ct^c) + 1/(Ct^C) < 1 - 1/(Ct^C)$. Finally, let $j \in [t]$ and $a \in \Sigma$ be the indices that minimize the statistical distance on the left-hand side of (8). Since $f(X_e|_{e_{ja} \notin X_e}) = A(X|_{X_j \neq a})$ and $f(X_e|_{e_{ja} \in X_e}) = A(X|_{X_j = a})$, we obtain the claim with $v \doteq e_{ja}$. \square

4.2 Proof of Theorem 4

Proof (of Theorem 4) Let A be the assumed algorithm compressing $t(n)$ instances of length n each to $\Gamma \cdot t \log t$ bits, where Γ is some positive constant. As in the proof of Theorem 3, we use A to construct a deterministic P/poly reduction from L to a conjunction of polynomially many instances of the statistical distance problem $\text{SD}_{\leq \delta}^{\geq \Delta}$.

This time we set $\Delta = 1$ and $\delta = 1 - \frac{1}{\text{poly}(n)}$; we will be more precise about δ below.

Since there is a polynomial gap between δ and Δ , Theorem 2 implies that $\text{SD}_{\leq \delta}^{\geq \Delta}$ is contained in the intersection of NP/poly and coNP/poly. Since the intersection is closed under polynomial conjunctions, we obtain $L \in \text{NP/poly} \cap \text{coNP/poly}$. Thus it remains to find such a reduction.

Lemma 8 applied to A yields a constant $C > 0$ so that the following hypergraph tournament $S: \binom{V}{Ct^C \cdot t} \rightarrow V$ with $V = \bar{L}_n$ is well-defined:

$$S(e) \doteq \min \left\{ v \in e \mid d \left(A(X_e|_{v \notin X_e}), A(X_e|_{v \in X_e}) \right) \leq \delta \right\}.$$

Lemma 8 actually requires some partition $e = e_1 \dot{\cup} \dots \dot{\cup} e_t$ with $|e_1| = \dots = |e_t| = Ct^C$ to be chosen, and we assume this to be done in the canonical lexicographic fashion. As already mentioned, we set $\Delta = 1$, and we set $\delta = 1 - 1/(Ct^C)$ to be equal to the right-hand side of (9).

To construct the advice of the reduction, note that, if $|V| \leq Ct^C \cdot t = \text{poly}(n)$, then we can directly use V as the advice. Otherwise, the advice at input length n is the dominating set $\mathcal{D} \subseteq \binom{V}{Ct^C \cdot t-1}$ whose existence is guaranteed by Lemma 5; in particular, its size is bounded by $t \cdot Ct^C \cdot n = \text{poly}(n)$.

The algorithm for L that uses $\text{SD}_{\leq \delta}^{\geq \Delta}$ as an oracle works as follows: Let $v \in \{0, 1\}^n$ be an instance of L given as input. If $v \in g$ holds for some $g \in \mathcal{D}$, the reduction rejects v and halts. Otherwise, for each $g \in \mathcal{D}$, it queries the SD-oracle on the instance $(A(X_e|_{v \notin X_e}), A(X_e|_{v \in X_e}))$ with $e = g \cup \{v\}$. If the oracle claims that all queries are yes-instances, our reduction accepts, and otherwise, it rejects.

The correctness of this reduction is analogous to the proof Theorem 3: If $v \in L$, then Lemma 2 guarantees that the statistical distance of all queries is one, and so all queries will detect this. If $v \in \bar{L}$, then since \mathcal{D} is a dominating set of S , we have $v \in g$ or $S(g \cup \{v\}) = v$ for some $g \in \mathcal{D}$. The latter will be detected in the query corresponding to g since $\delta < \Delta$. This completes the proof of the theorem. \square

From its proof, we observe that Theorem 4 would even apply to randomized OR-compressions A that have a two-sided error of up to $1/\text{poly}(n)$, so long as the polynomial in the denominator is large enough so that the gap between Δ and δ remains at least an inverse polynomial of n . Handling larger error probabilities does not seem possible with the technique that relies on Vajda's inequality.

4.3 Information-Theoretic Arguments

Proof (of Lemma 7) Vajda's inequality [20,21] is an alternative way to bound the statistical distance in terms of the KL-divergence. We use the following version of the inequality (cf. [9, Theorem 4.8]):

$$d(X, Y) \leq 1 - \frac{1}{e \cdot 2^{D_{\text{KL}}(X \| Y)}}.$$

We use the bound as follows:

$$\begin{aligned} \mathbf{E}_j \mathbf{E}_a d(f(X), f(X|_{j \leftarrow a})) \\ \leq \mathbf{E}_j \mathbf{E}_a \left(1 - \exp \left(-1 - D_{\text{KL}}(f(X|_{j \leftarrow a}) \| f(X)) \right) \right) \end{aligned} \quad (\text{Vajda's inequality})$$

$$\begin{aligned} &\leq 1 - \exp\left(-1 - \mathbf{E}_j \mathbf{E}_a D_{\text{KL}}\left(f(X|_{j \leftarrow a}) \parallel f(X)\right)\right) && \text{(Jensen's inequality)} \\ &\leq 1 - e^{-1 - \Gamma \log t} \leq 1 - 1/(ct^c). && \text{(Lemma 6)} \end{aligned}$$

Here, setting $c = O(\Gamma)$ is sufficient. Now (8) follows from the triangle inequality:

$$\begin{aligned} \mathbf{E}_j \mathbf{E}_a d\left(f(X|_{X_j \neq a}), f(X|_{X_j = a})\right) &\leq \mathbf{E}_j \mathbf{E}_a d\left(f(X|_{X_j \neq a}), f(X)\right) \\ &\quad + \mathbf{E}_j \mathbf{E}_a d\left(f(X), f(X|_{X_j = a})\right) \\ &\leq \frac{1}{|\Sigma|} + 1 - \frac{1}{ct^c}. \end{aligned}$$

For this, we use the following basic inequalities of the statistical distance as defined in (1):

$$\begin{aligned} d(f(X|_{X_j \neq a}), f(X)) &\leq d(X|_{X_j \neq a}, X) \\ &\leq \Pr(X_j \neq a) \cdot d(X|_{X_j \neq a}, X|_{X_j \neq a}) + \Pr(X_j = a) \cdot d(X|_{X_j \neq a}, X|_{X_j = a}) \\ &\leq \Pr(X_j \neq a) \cdot 0 + \Pr(X_j = a) \cdot 1 = \Pr(X_j = a). \end{aligned}$$

The latter equals $\frac{1}{|\Sigma|}$ since X_j is uniformly distributed on Σ . □

5 Extension: f -Compression

We end this paper with a small observation: Instead of OR-compressions or AND-compressions, we could just as well consider f -compressions for a Boolean function $f: \{0, 1\}^t \rightarrow \{0, 1\}$. If the function f is symmetric, that is, if $f(x)$ depends only on the Hamming weight of x , then we can represent f as a function $f: \{0, \dots, t\} \rightarrow \{0, 1\}$. We make the observation that Drucker's theorem applies to f -compressions whenever f is a non-constant, symmetric function.

Definition 3 Let $f: \{0, \dots, t\} \rightarrow \{0, 1\}$ be any function. Then an f -compression of L into L' is a mapping

$$A: \binom{\{0, 1\}^n}{\leq t} \rightarrow \{0, 1\}^{\epsilon t},$$

such that, for all $x \in \binom{\{0, 1\}^n}{\leq t}$, we have $A(x) \in L'$ if and only if $f(|x \cap L|) = 1$.

We list some examples:

- OR-compressions are f -compressions with $f(i) = 1$ if and only if $i > 0$.
- AND-compressions are f -compressions with $f(i) = 1$ if and only if $i = t$.
- Majority-compressions are f -compressions with $f(i) = 1$ if and only if $i > t/2$.
- Parity-compressions are f -compressions with $f(i) = 1$ if and only if i is odd.

We next prove that Theorem 3 and 4 can be applied whenever f is a symmetric function that is not constant. Using more complicated arguments, Drucker [9] already achieves this for a more general class of Boolean functions f .

Lemma 9 *Let $f: \{0, \dots, t\} \rightarrow \{0, 1\}$ be not constant. Then every f -compression for L with size ϵt can be transformed into a compression for L or for \bar{L} , in the sense of Theorem 3 and with size bound at most $2\epsilon t$.*

Proof Let A be an f -compression from L into L' . Then A is also a $(i \mapsto 1 - f(i))$ -compression from L into \bar{L}' , an $(i \mapsto f(t - i))$ -compression from \bar{L} into L' , and an $(i \mapsto 1 - f(t - i))$ -compression from \bar{L} into \bar{L}' . Since f is not constant, at least one of these four views corresponds to a function f' for which there is an index $i \leq t/2$ so that $f'(i) = 0$ and $f'(i + 1) = 1$ hold. Assume without loss of generality that this holds already for f . Then we define $A': \binom{[0, 1]^n}{\leq t-i} \rightarrow \{0, 1\}^{\epsilon t}$ as follows:

$$A'(\{x_{i+1}, x_{i+1}, \dots, x_t\}) \doteq A(\{\top_1, \dots, \top_i, x_{i+1}, x_{i+1}, \dots, x_t\}),$$

where \top_1, \dots, \top_i are arbitrary distinct yes-instances of L at input length n . For the purposes of Theorem 3, these instances can be written in the non-uniform advice of A' . If this many distinct yes-instances do not exist, then the language L is trivial at this input length. To ensure that the x_j 's are distinct from the \top_j 's, we actually store a list of $2t$ yes-instances \top_j and inject only i of those that are different from the x_j 's.

A' is just like A , except that i inputs have already been fixed to yes-instances. Then A' is a compressive map that satisfies the following: If $|x \cap L| = 0$ then $A'(x) \notin L'$, and if $|x \cap L| = 1$ then $A'(x) \in L'$. Since the number of inputs has decreased to $t' = t - i \geq t/2$, the new size of the compression is $\epsilon t \leq 2\epsilon t'$ in terms of t' . \square

Acknowledgments I would like to thank Andrew Drucker, Martin Grohe, and Noy Galil Rotbart for encouraging me to pursue the publication of this manuscript, David Xiao for pointing out Theorem 2 to me, Andrew Drucker, Dániel Marx, and anonymous referees for comments on an earlier version of this paper, and Dieter van Melkebeek for some helpful discussions.

References

1. Drucker, A.: New limits to classical and quantum instance compression. In: Proceedings of the 53rd Annual Symposium on Foundations of Computer Science (FOCS), pp. 609–618 (2012). doi:[10.1109/FOCS.2012.71](https://doi.org/10.1109/FOCS.2012.71)
2. Bodlaender, H.L., Downey, R.G., Fellows, M.R., Hermelin, D.: On problems without polynomial kernels. J. Comput. Syst. Sci. **75**, 423–434 (2009)
3. Fortnow, L., Santhanam, R.: Infeasibility of instance compression and succinct PCPs for NP. J. Comput. Syst. Sci. **77**, 91–106 (2011)
4. Ko, K.-I.: On self-reducibility and weak P-selectivity. J. Comput. Syst. Sci. **26**, 209–211 (1983)
5. Harnik, D., Naor, M.: On the compressibility of NP instances and cryptographic applications. SIAM J. Comput. **39**, 1667–1713 (2010)
6. Dell, H., van Melkebeek, D.: Satisfiability allows no nontrivial sparsification unless the polynomial-time hierarchy collapses. J. ACM **61** (2014). doi:[10.1145/2629620](https://doi.org/10.1145/2629620)
7. Dell, H., Kabanets, V., van Melkebeek, D., Watanabe, O.: Is Valiant-Vazirani's isolation probability improvable? Comput. Complex. **22**, 345–383 (2013)
8. Xiao, D.: New perspectives on the complexity of computational learning, and other problems in theoretical computer science Ph.D. thesis. Princeton University. [ftp://ftp.cs.princeton.edu/techreports/2009/866.pdf](http://ftp.cs.princeton.edu/techreports/2009/866.pdf) (2009)

9. Drucker, A.: New limits to classical and quantum instance compression tech report TR12-112 rev. 3 (Electronic Colloquium on Computational Complexity (ECCC). <http://eccc.hpi-web.de/report/2012/112/> (2014)
10. Dell, H., Marx, D.: Kernelization of packing problems. In: Proceedings of the 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 68–81 (2012). doi:[10.1137/1.9781611973099.6](https://doi.org/10.1137/1.9781611973099.6)
11. Hermelin, D., Wu, X.: Weak compositions and their applications to polynomial lower bounds for kernelization. In: Proceedings of the 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 104–113 (2012). doi:[10.1137/1.9781611973099.9](https://doi.org/10.1137/1.9781611973099.9)
12. Bodlaender, H.L., Jansen, B.M.P., Kratsch, S.: Kernelization lower bounds by cross-composition. *SIAM J. Discrete Math.* **28**, 277–305 (2014)
13. Kratsch, S.: Co-nondeterminism in compositions: a kernelization lower bound for a Ramsey-type problem. In: Proceedings of the 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 114–122 (2012). doi:[10.1137/1.9781611973099.10](https://doi.org/10.1137/1.9781611973099.10)
14. Kratsch, S., Philip, G., Ray, S.: Point line cover: the easy kernel is essentially tight. In: Proceedings of the 25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 1596–1606 (2014). doi:[10.1137/1.9781611973402.116](https://doi.org/10.1137/1.9781611973402.116)
15. Arora, S., Barak, B.: Computational Complexity—A Modern Approach, Cambridge University Press. ISBN: 978-0-521-42426-4 (2009)
16. Adleman, L.M.: Two theorems on random polynomial time. In: Proceedings of the 19th Annual Symposium on Foundations of Computer Science (FOCS), pp. 75–83 (1978). doi:[10.1109/SFCS.1978.37](https://doi.org/10.1109/SFCS.1978.37)
17. Sahai, A., Vadhan, S.: A complete problem for statistical zero knowledge. *J. ACM* **50**, 196–249 (2003)
18. Goldreich, O., Vadhan, S.: On the complexity of computational problems regarding distributions (a survey) Tech report TR11-004 Electronic Colloquium on Computational Complexity (ECCC). <http://eccc.hpi-web.de/report/2011/004/> (2011)
19. Cover, T.M., Thomas, J.A.: Elements of Information Theory. Wiley, London (2012)
20. Fedotov, A.A., Harremoës, P., Topsøe, F.: Refinements of Pinsker’s inequality. *IEEE Trans. Inf. Theory* **49**, 1491–1498 (2003)
21. Reid, M., Williamson, B.: Generalised Pinsker inequalities. In: Proceedings of the 22nd Annual Conference on Learning Theory (COLT), pp. 18–21. <http://www.cs.mcgill.ca/~colt2009/papers/013.pdf> (2009)